

# Learning based Clustering for the Automatic Annotations from Web Databases

Richa Saxena  
M.Tech CSE  
SRCEM, Banmore

Sushil Kumar Chaturvedi  
Asst. Prof.  
SRCEM, Banmore

## ABSTRACT

Rapid increase of use of internet provides knowledge extraction from the web databases and HTML pages associated with it. Although there are various techniques implemented for the access of the annotations of the search results from the web databases. Here in this paper by identifying the problems with the existing techniques for the annotation search results from web databases such as alignment problem or to split composite text node when there are no explicit separators.

Here propose an efficient technique which overcomes the above problems by using some supervised learning algorithm such as support vector machine. The technique implemented provides high rate of information by providing high annotations search results from web databases.

The proposed method implemented here for the efficient retrieval of text nodes and data units using supervised learning approach using SVM provides efficient precision and recall as compared to the existing approach.

The proposed methodology implemented here using SVM based clustering and labeling of search records is compared with existing methodology implemented for the search records. The Result Analysis shows the performance of the proposed methodology.

The proposed method shows higher precision and recall as well as has high Accuracy for the prediction of annotated search records from the web databases.

## Keywords

Annotations, Wrapper, Semantic Model, HTML Tags, NLP, Ontology, UIUC.

## 1. INTRODUCTION

Formerly with the growth of web-based resources, including an explosion of user-generated content, has come parallel growth of research into web-based searching behaviour and searcher experiences. Also simultaneous with advancements in web technologies and relevance ranking algorithms, search engines have achieved a towering level of trust and a perception of the search engine's mysterious inner workings that verges on clairvoyance: simply type in a stream of consciousness and the oracle-like search engine responds with a useful answer. After pressing search button, scan the results list, and select a few from the first page or two of results. However, this simplistic perspective may be contributing to a lack of understanding of the information environment, leaving students in a world of possible impediments to searching, without an understanding of ways to improve the process. The Web has become the preferred medium for many database forms and usages used to store information. Database-driven Web sites have their own interfaces and access forms for creating HTML pages on the fly. Web database technologies

define the way that these forms can connect to and retrieve data from database servers.[1] The number of database-driven Websites is increasing and they do not send queries to Web databases. World Wide Web has resulted in a huge amount of information sources on the Internet. Web information sources, access to this huge collection of information has been limited to browsing and searching.

Web mining seems to be the part of data mining that is used for the extraction of information from the knowledge databases. Web mining is used for the extraction of web contents and hence works on the basis of Web content and web usage data based and web structure based.

Web content mining is the process of extracting useful information from the contents of web documents and texts. Hence on the basis of Content of the data is the collection of a web page which is designed to hold. These data contents consist of various text and images as well as videos. Research activities on this topic have drawn heavily on techniques developed in other disciplines such as Information Retrieval (IR) and Natural Language Processing (NLP). While there exists a significant body of work in extracting knowledge from images in the fields of image processing and computer vision.

Web usage mining is the application of data mining techniques to discover interesting usage patterns from web usage data, in order to understand and better serve the needs of web-based applications.

Efficiency of searching and updating information increases by Alignment and annotation of statistics. Alignment of data can be referred to as arranging the data in such a way that data inside the same group have the same meaning and accessing. It is a methodology for adding information to a text sequences such as article. Data annotation [2] enables fast retrieval of information in the deep web. A data unit is a part of text that semantically represents real world entity concepts. Dynamically for human browsing these data units are encoded into the result page and assigned meaningful labels. Annotate the data units requires lots of human efforts.

### 1.1 Annotation Phases

*Phase 1:* Alignment phase: In alignment phase align all the data into different group. Here most of the used group belongs to a different concept.

*Phase 2:* Annotation phase: In annotation phase used several basic annotators with each exploiting one type of features.

*Phase3:* Annotation wrapper generation phase: In annotation wrapper generation phase an annotation rule is generated for each identified entity or concept.

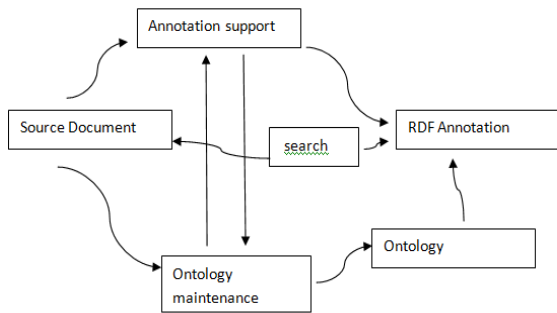


Figure 1: Phases of automatic annotation solution [23]

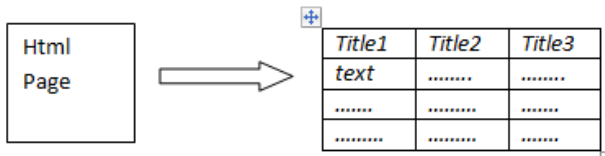


Figure 2: Extracts (automatically) text from a web-page into a table [2]

## 1.2 Data unit and text node relationships

Data unit [3] is a piece of text that semantically represents concept of real world entity. Data unit is totally different from text node is visible element on the web page and data unit located in the text nodes.

*One-to-One Relationship:* (referred as atomic text nodes). Text node involves only one data unit

*One-to-Many Relationship:* (referred as composite text nodes) a text node consists of multiple data units i.e. multiple data units are encoded into single text nodes.

*Many-to-One Relationship:* (referred as decorative tags) multiple text nodes are encoded into single data unit. This type of text nodes is referred as decorative tags because they are used for changing the appearance of part of the text node.

*One-To-Nothing Relationship:* (referred as template text nodes) Text nodes are not part of any data unit inside SRRs. This relationship for text nodes and data units are represents the relation in between them. This type of text node is referred as template text node.

There are five common features shared by the data units Data content Presentation style Data type Tag path Adjacency

Data content [3]: Data unit or text node of same concept shares certain keywords which are used to search the information quickly. For e.g., keyword “machine” will return the information that are relevant to word machines.

Presentation style: Presentation feature describes how a data unit is displayed on a web.

Data alignment and labeling [3]: Current tasks when compared with automatic annotation approach. They are based on one or a few facilities. Automatic annotation alignment approach first data units and handles relations between text nodes and data unit do use variety of features the device is a cluster-based transfer algorithm and is used in the alignment process. Label assignment IIS (unified interface schema) and LIS (local interface schema). There are attributes in all LIS IIS and thus eliminates inadequacy and inconsistent labels label problems. In the coalition of some basic

annotators groups started to annotate and combine multiple annotators a probability model is used for the results of this approach are called multiple-annotator approach.

Consider a set of SRRs that are extracted from a result page returned from the web database. The Automatic annotation approach has three major phases as shown in the [4].

Alignment phase [3]: First data alignment phase in the SRRs units identified and organized into different groups for each group corresponds to a different concept (for example, all titles of books are grouped together). Figure 1b across all SRRs step 1 each column containing data unit with same sense results. This step is to identify the patterns and features of data between units are used.

## 2. RELATED WORK

In this paper author has [5] states that there are a growing number of pages that can't be indexed by search engines and stay invisible to other surfers, despite the fact that they contain a lot of relevant content. Pages with dynamic content i.e the pages that are the result of a submitted query and consequently do not have a static URL can be impossible for search engines to find, since the crawlers cannot replicate the query submission carried out by human beings. Generally traditional search engines follow links to the index page on a site and then crawl from there to other pages by following links. Search engine crawlers will therefore have more difficulties seeing a page that is not linked to from any other page(hidden web pages).

Annotation [4] a collaborative client server system document annotation is a special they are stored on the server in such a way that anyone who has access to an annotation server for a given document to consult all related annotation and add your own annotations will be enabled for these annotations are divided into typing comments Improve projections, assumptions. This system was developed using W3C standards. Yet, only possible Committee on State annotation text; It is annotated by a picture or symbol. EXCOM [4] is an annotation engine internal/external annotating a document by a clique of hay knowledge on aim uses a set of linguistic devices. This engine is under development, and at the present time, a cosmic stories and questions allows the production of an expressed, since this technology is not entirely

Here in this paper author has proposed a new method [6] Acacia team allows annotation system developed by genes. This creature, which experiments to validate and to interpret the results, obtained on the biopuces helps make system research difficult task them. Its genetic database offers the possibility of a key word research. For key word can or a biological phenomenon Jean correspondence study. All previous works are interested in general documents annotation like scientific articles, Web documents, biological databases and multimedia documents. Only few of them focus on the events annotation. Here in this paper they present, in the following, some of these works: The annotation of temporal information in texts [6]: this work focused more specifically on relations between events introduced by verbs in finite clauses. It proposes a procedure that achieves the task of annotation and a way of measuring the results. The authors of this work tested the feasibility of this procedure on newswire articles with promising results. Then, they developed two evaluation measures of the annotation: fineness and consistency.

Annotating texts [7] features and relationships to determine the relative annotation scheme: This enables order and, if possible, absolute time events. A planning an annotated corpus can be used for building the corpus is usually producing such benefits associated with building resources. It also can be used to better understand the phenomena. Plus it training and adaptive algorithms for evaluating represents a source. It automatically shows the relationship of the features and interest. However, we observed that the relationship between the incidences of this work to determine based on temporal markers only. There are inherent differences with regard to events without using temporal markers which are accurate.

SyDoM [8] is a semantic annotation of Web pages system. This allows the enrichment of these pages so that they take account of their writing without language find it with textual XML format is dedicated to manage documents stored. we see that SyDoM has two main advantages: first, multilingual research and other But the improvement of the representation of Web pages, we SyDoM [8] out research on Web pages only if it has been already annotated, yet this annotation by using different means to inquire Web pages created thesauri have been unable to get that information.

The W3C also made the task of making existing databases available for the Semantic Web one of its goals and initiated the RDB2RDF Incubator Group. In the course of their work in 2009, they collected and evaluated the state of the art in this field and published their final report in [9]. This survey showed several approaches – from a complete transformation of an existing relational database to an RDF database on one side, to on-demand mapping and query translations from SPARQL to SQL on the other side. When publishing data that has to be maintained and updated over time, it is impractical to have to publish the same information twice – once for humans and a second time for tools, especially if content is also provided by users of the site. In this case, it is not just a single effort when initially publishing the data, but results in consequently having to update and maintain two separate pieces of work. A practical and easy way to integrate semantic information into an existing document is to use annotations. A first solution for this was proposed by [9] who described the concept of miro formats. These are small sets of semantic data that can be embedded in a webpage, invisible to the user but visible for tools and search engines. By using this data, structured information about things like authorship or even cooking recipes can be given that can be extracted from the page. But these formats have to be agreed on by the community in order to understand the structure and the content.

Only in 2008 Google has started to develop strategies to “surface” information from this vast source: in Ref. [10], employees of Google describe their solution to automatically extract information from databases by creating queries with sets of keywords depending on the apparent topic of the site. But this method can only extract parts of the data in the database. According to Ref. [10] the success depends mainly on the size of the database and was as low as 20% in the documented results.

We believe that these web-databases could be a great opportunity to show the strengths of the Semantic Web. If a practical solution can be found to make these databases usable for semantic agents, there would be no need to index all of its content, because it could be accessed on demand. And the prospect of having 400-550 times more information available than current search engines can access would instantly solve

the chicken-and-egg problem and make it worthwhile to develop tools and services for the Semantic Web.

LingLiu XWRAP is based on the concept of XML-enabled wrapper for various sources. The paper [11] describes the methodology and the software development of XWRAP.

Crescenzi,V [12] Efficient Techniques for Effective Wrapper Induction several studies have recently concentrated on the generation of wrappers for extracting data from Web data sources.

### **3. PROPOSED METHODOLOGY**

The alignment algorithm implemented here is for composite text nodes and other tags available in the HTML tags.

1. Merge Text nodes: Here detection and removal of the decorative tags can be identified to merge them into single text nodes or into multiple text nodes.
2. Align Text nodes: Here the HTML tags which contain some meaningful labels can be aligned so that it can be used for efficient searching.
3. Split Text nodes: The multiple text nodes or the HTML tags containing composite nodes can be split here to provides and generate multiple search results.
4. Align Data units: The step is carried out for the separation of composite groups into multiple aligned groups.

The proposed methodology implemented here uses the following Alignment algorithm for the single or multiple text nodes and then uses SVM based clustering to cluster the similar text nodes which provide same set of search results from the HTML tags.

The proposed algorithm that is implemented here consists of the following basis steps:

#### **Step 1**

Take an input dataset which contains a number of web pages for various categories such as Book and Movies and Pen Drives and Electronics, since these categories contains some basis labels such as title and author and price and ISBN for the category Book. The Web pages taken here contains HTML tags and nodes inside which some meaning information is stored.

#### **Step 2**

As soon as the input dataset is selected the next step is to extract the important features from the web pages. Here for the extraction of features for the web pages the following predefined classes are used which removes the decorative tags from the web pages and stores text nodes.

```
ParserDelegator delegator = new ParserDelegator();
```

```
delegator.parse(in, this, Boolean.TRUE);
```

#### **Step 3**

After the extraction of features from each of the input web page dataset cosine similarity is measure for each of the text document web pages using the following formula?

$$\text{VectAB} = \text{VectAB} + (\text{freq1} * \text{freq2});$$
$$\text{VectA\_Sq} = \text{VectA\_Sq} + \text{freq1} * \text{freq1};$$
$$\text{VectB\_Sq} = \text{VectB\_Sq} + \text{freq2} * \text{freq2};$$

$$\text{sim\_score} = \frac{V_{AB}}{\sqrt{V_A} * \sqrt{V_B}}$$

$$V_{AB} = V_{AB} + (F_A * F_B)$$

Where,

$$V_{AB} = V_{AB} + (F_A * F_B)$$

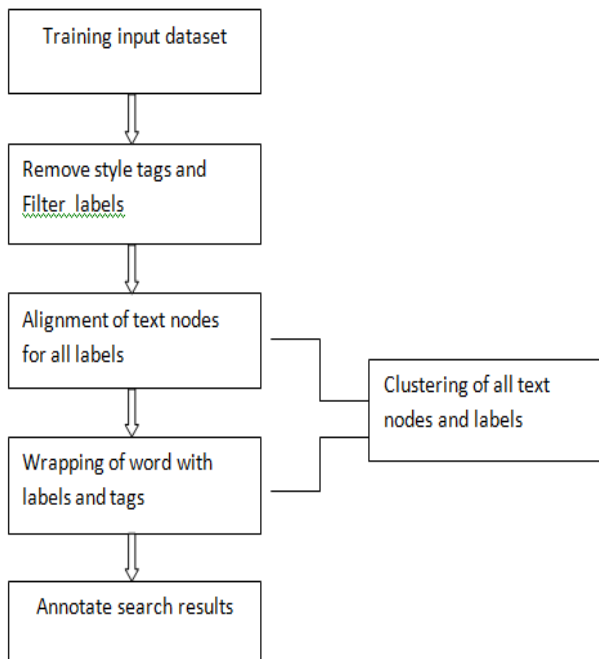
**Step 4**

The next step is the computation of Alignment of the text nodes and the data units that are available in the web pages. Here Alignment can be done using the clustering using Support vector machine. A Support vector machine contains two types of process one is linear based and other is Gaussian based also known as radial basis function. SVM is an efficient learning approach which takes ‘X’ as an input data values and ‘Y’ as the label values and gamma coefficient and Class index ‘C’.

**Step 5**

Finally on the basis of clustering of the labels of the text nodes value a weighted factor is decided and hence labels are assigned.

The figure in 3 shown below is the flow diagram of the proposed methodology. The method implemented here consists of a selection of a number of input web pages. Since these web pages contain HTML tags hence the next step is the removal of these decorative tag nodes and hence text nodes are extracted. These Text nodes contain a number of data units which contains some meaningful information which may be one or many. The data units are then aligned using supervised learning approach such as support vector machine which provides a number of assigned labels to these data units.



**Figure 3 Flow Diagram of the Proposed Methodology**

**4. RESULT ANALYSIS**

The table shown below in 1 is the experimental analysis of the searching of annotations using label based clustering. The result is analysed on various domains such as Book and Pen Drives and Music and Movies as well as Games. The result is analysed on the basis of Precision and Recall and F-Score. Here Precision can be computed on the basis of correctly identified annotations to the total number of annotations fetched from the web databases. Recall is the computation of total number of annotations fetch from the web databases to the total number of annotated records present.

**Table 1 Analysis of Existing work on various Domains**

Domain	Precision	Recall	F-Score
Book	0.5	0.5	0.5
Pen Drive	0.368	0.632	0.465
Music	0.4	0.6	0.48
Movies	0.46	0.56	0.505098
Games	0.58	0.47	0.519238

The table shown below in 2 is the experimental analysis of the searching of annotations using SVM based clustering. The result is analysed on various domains such as Book and Pen Drives and Music and Movies as well as Games. The result is analysed on the basis of Precision and Recall and F-Score. Here Precision can be computed on the basis of correctly identified annotations to the total number of annotations fetched from the web databases. Recall is the computation of total number of annotations fetch from the web databases to the total number of annotated records present.

**Table 2 Analysis of Proposed work on various Domains**

Domain	Precision	Recall	F-Score
Book	0.65	0.54	0.59
Pen Drive	0.825	0.548	0.66
Music	0.754	0.712	0.7324
Movies	0.823	0.743	0.781
Games	0.632	0.593	0.6119

The figure shown below in 4 is the experimental analysis and comparison of Precision based on Existing and Proposed work. The result is analysed on various domains such as Book and Pen Drives and Music and Movies as well as Games. Here Precision can be computed on the basis of correctly identified annotations to the total number of annotations fetched from the web databases.

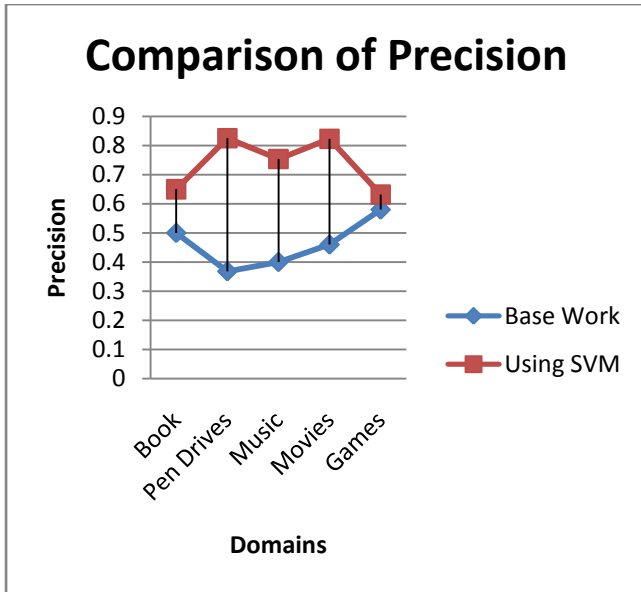


Figure 4 Comparison of Precision

The figure shown below in 5 is the experimental analysis and comparison of Recall based on Existing and Proposed work. The result is analysed on various domains such as Book and Pen Drives and Music and Movies as well as Games. Recall is the computation of total number of annotations fetch from the web databases to the total number of annotated records present.

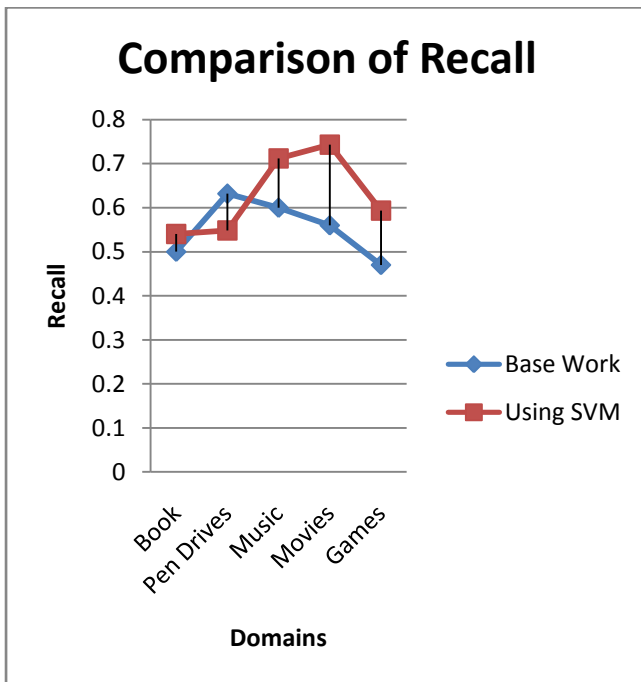


Figure 5 Comparison of Recall

The figure shown below in 6 is the experimental analysis and comparison of F-Score based on Existing and Proposed work. The result is analysed on various domains such as Book and Pen Drives and Music and Movies as well as Games. It is defined as:

$$F - Score = \frac{2 * precision * recall}{precision + recall}$$

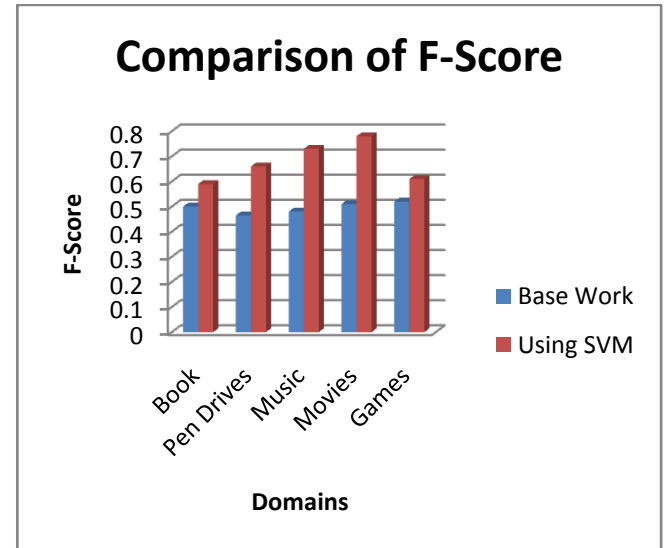


Figure 6 Comparison of F-Score

## 5. CONCLUSION AND FUTURE WORK

The comparison between existing and proposed techniques can be analyzed based on precision and recall and it was found that the proposed methodology not only removes the problems of the existing technique but also provides high precision and recall. The proposed methodology implemented here for the searching of the annotations from the web databases. Here the annotations can be identified on the various categories such as Book, Movies, Electronics, Pen Drives, Auto. The proposed methodology is applied on these categories with different web pages and hence on the basis of search web records labels are assigned to these web pages. After identification of annotations in the web databases accuracy can be computed and compared to the existing technique that is implemented for the efficient search of records from the web databases and the proposed methodology provides high precision and recall as compared to the existing technique.

Although the technique is efficient and provides efficient results as compared to the existing technique but further enhancements can be done for the improvement of accuracy for other keywords and annotations.

## 6. REFERENCES

- [1] Y. Lu, H. He, H. Zhao, W. Meng, and C. Yu, "Annotating Structured Data of the Deep Web," Proc. IEEE 23rd Int'l Conf. Data Eng. (ICDE), 2007.
- [2] Priyanka P. Boraste "A Survey on Data Annotation for the Web Databases" IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661, p- ISSN: 2278-8727 Volume 16, Issue 2, Ver. XI (Mar-Apr. 2014), PP 68-70 www.iosrjournals.org.
- [3] Y. Pauline Jeba, Mrs. P. Rebecca Sandra, "A Survey On Annotating Search Results From Web Databases", International Journal Of Research In Computer Applications And Robotics, Vol -1, Issue-9, 2013.
- [4] J. Kahan, M-R. Koivunen, Annotea: an open RDF infrastructure for shared Web annotations. Proceedings of the 10<sup>th</sup> international conference on World Wide Web, 2001.

- [5] L. Gravano, H. Garcia-Molina, A. Tomasic, "GLOSS: Text-Source Discovery over Internet", *TODS* 24(2), 1999.
- [6] K. Khelif, R. Dieng-Kuntz, P. Barbry, An Ontology-based Approach to Support Text Mining and Information Retrieval in the Bio logical Domain, in *J. UCS* 13(12), pp. 1881-1907, 2007.
- [7] A. Setzer, R. Gaizauskas, TimeM L: Robust specification of event and temporal expressions in text. In *The second international conference on language resources and evaluation*, 2000.
- [8] C. Roussey, S. Calabretto, An experiment using Conceptual Graph Structure for a Multilingual Information System, in the 13<sup>th</sup> International Conference on Conceptual Structures, ICCS'2005.
- [9] A Survey of Current Approaches for Mapping of Relational Databases to RDF. Retrieved October 28, 2011 from [www.w3.org/2005/Incubator/rdb2rdf/RDB2RDF\\_SurveyReport.pdf](http://www.w3.org/2005/Incubator/rdb2rdf/RDB2RDF_SurveyReport.pdf), 2005.
- [10] J. Madhayan et al, "Google's Deep-Web Crawl." *Proceedings of the VLDB Endowment*, Vol. 1, Issue 2, pp. 1241-1252, 2008.
- [11] A Survey of Web Information Extraction Systems Chia-Hui Chang, Member, IEEE Computer Society, Mohammed Kayed, Moheb Ramzy Girgis, Member, Ieee Transactions On Knowledge And Data Engineering, VOL. 18, NO. 10, OCTOBER 2006
- [12] V. Crescenzi, G. Mecca, and P. Merialdo, "RoadRUNNER: Towards Automatic Data Extraction from Large Web Sites," *Proc. Very Large Data Bases (VLDB) Conf.*, 2001.